

WHAT IS CLAIMED IS:

1. A storage comprising:

5 a non-volatile memory storing a first inode corresponding to a first file; and

 a block manager configured to copy said first inode to a second inode, wherein
 said block manager is configured to change said second inode in response
 to updates to said first file, and wherein said block manager is configured
10 to atomically update said first file in response to a commit of said first file
 by writing said second inode to said non-volatile memory.

2. The storage as recited in claim 1 wherein said non-volatile memory stores a journal
comprising a list of committed inodes, and wherein said block manager is configured to
15 record said second inode in said journal.

3. The storage as recited in claim 2 wherein said commit of said first file comprises a
commit command received from an external source which updates said first file.

20 4. The storage as recited in claim 3 wherein said commit command comprises a file close
command.

5. The storage as recited in claim 3 wherein said commit command comprises an fsync
command.

25 6. The storage as recited in claim 2 wherein said journal further includes a checkpoint
record including a description of an inode file, a block allocation bitmap, and an inode
allocation bitmap.

7. The storage as recited in claim 6 wherein the description comprises inodes for each of said inode file, said block allocation bitmap, and said inode allocation bitmap.

8. An apparatus comprising:

5

a computing node configured to perform one or more write commands to a file
and a commit command committing the one or more write commands to
said file; and

10

a storage coupled to receive said one or more write commands and said commit
command, wherein said storage is configured to copy one or more blocks
of said file to a copied one or more blocks, said one or more blocks
updated by said one or more write commands, and wherein said storage is
configured to update said copied one or more blocks with write data
15 corresponding to said one or more write commands, and wherein said
storage is configured to copy a first inode corresponding to said file to a
second inode and to update pointers within said second inode
corresponding to said one or more blocks to point to said copied one or
more blocks, and wherein said storage is configured to atomically update
20 said file by writing said second inode responsive to said commit
command, and wherein said first inode is stored in an inode file, and
wherein said inode file is identified by a master inode, and wherein said
inode file is atomically updated with said second inode by writing said
master inode subsequent to said commit command.

25

9. The apparatus as recited in claim 6 wherein said commit command comprises a file
close command.

10. The apparatus as recited in claim 6 wherein said commit command comprises an

fsync command.

11. A method comprising:

- 5 copying a first inode corresponding to a first file to a second inode;
- modifying said second inode in response to one or more changes to said first file;
- and
- 10 atomically updating said first file by establishing said second inode as the inode
 for said first file.

12. The method as recited in claim 11 wherein said establishing comprises storing said
second inode in a journal stored in a nonvolatile memory.

- 15 13. The method as recited in claim 12 further comprising writing a master inode
 corresponding to an inode file including said second inode to a checkpoint record in said
 journal.

- 20 14. The method as recited in claim 13 wherein recovering from a system failure
 comprises:

 scanning said journal to locate a most recent checkpoint record and zero or more
 inodes subsequent to said most recent checkpoint record within said
25 journal;

 copying said master inode from said most recent checkpoint record to a volatile
 memory; and

updating an inode file corresponding to said master inode with said one or more inodes subsequent to said most recent checkpoint record.

15. The method as recited in claim 14 wherein said updating said inode file comprises:

5

copying one or more blocks of said inode file storing said one or more inodes to a copied one or more blocks; and

10

updating said master inode in said volatile memory to point to said copied one or more blocks.

16. The method as recited in claim 11 wherein said block map further comprises a first inode allocation bitmap indicating which inodes within said first inode file are allocated to files, the method further comprising:

15

copying said first inode allocation bitmap to a second inode allocation bitmap;

modifying said second inode allocation bitmap to reflect one or more inodes allocated to new files; and

20

establishing a second inode within said block map to said second inode allocation bitmap subsequent to said modifying said second inode bitmap.

17. The method as recited in claim 16 wherein said block map further comprises a first block allocation bitmap indicating which blocks within a storage including said block map are allocated to files, the method further comprising:

25

copying said first block allocation bitmap to a second block allocation bitmap;

modifying said second block allocation bitmap to reflect one or more blocks
allocated to files; and

establishing a third inode within said block map to said second block allocation
5 bitmap subsequent to said modifying said second block allocation bitmap.

18. The method as recited in claim 11 wherein said establishing said second inode is
performed in response to a commit command.

10 19. The method as recited in claim 18 wherein said commit command is a close file
command.

20. The method as recited in claim 18 wherein said commit command is an fsync
command.

15 21. A storage comprising:

a non-volatile memory storing a first inode corresponding to a first version of a
file; and

20 a block manager configured to copy said first inode to a second inode, wherein
said block manager is configured to change said second inode in response
to updates to the file, and wherein said block manager is configured to
atomically update the file, producing a second version of the file, in
25 response to a commit of the file by writing said second inode to said non-
volatile memory.

22. The storage as recited in claim 21 wherein said non-volatile memory stores a journal
comprising a list of committed inodes, and wherein said block manager is configured to

record said second inode in said journal.

23. The storage as recited in claim 22 wherein said commit of the file comprises a commit command received from an external source which updates the file.

5

24. The storage as recited in claim 23 wherein said commit command comprises a file close command.

25. The storage as recited in claim 23 wherein said commit command comprises an fsync
10 command.

26. The storage as recited in claim 22 wherein said journal further includes a checkpoint record including a description of an inode file, a block allocation bitmap, and an inode allocation bitmap.

15

27. The storage as recited in claim 26 wherein the description comprises inodes for each of said inode file, said block allocation bitmap, and said inode allocation bitmap.

28. A method comprising:

20

copying a first inode corresponding to a first version of a file to a second inode;

modifying said second inode in response to one or more changes to the file,
creating a second version of the file; and

25

atomically updating the file to the second version by establishing said second
inode as the inode for the file.

29. The method as recited in claim 28 wherein said establishing comprises storing said

second inode in a journal stored in a nonvolatile memory.

30. The method as recited in claim 29 further comprising writing a master inode
corresponding to an inode file including said second inode to a checkpoint record in said
5 journal.

31. The method as recited in claim 30 wherein recovering from a system failure
comprises:

10 scanning said journal to locate a most recent checkpoint record and zero or more
inodes subsequent to said most recent checkpoint record within said
journal;

15 copying said master inode from said most recent checkpoint record to a volatile
memory; and

updating an inode file corresponding to said master inode with said one or more
inodes subsequent to said most recent checkpoint record.

20 32. The method as recited in claim 31 wherein said updating said inode file comprises:

copying one or more blocks of said inode file storing said one or more inodes to a
copied one or more blocks; and

25 updating said master inode in said volatile memory to point to said copied one or
more blocks.

33. The method as recited in claim 28 wherein said block map further comprises a first
inode allocation bitmap indicating which inodes within said first inode file are allocated

to files, the method further comprising:

copying said first inode allocation bitmap to a second inode allocation bitmap;

5 modifying said second inode allocation bitmap to reflect one or more inodes
 allocated to new files; and

 establishing a second inode within said block map to said second inode allocation
 bitmap subsequent to said modifying said second inode bitmap.

10

34. The method as recited in claim 33 wherein said block map further comprises a first
block allocation bitmap indicating which blocks within a storage including said block
map are allocated to files, the method further comprising:

15 copying said first block allocation bitmap to a second block allocation bitmap;

 modifying said second block allocation bitmap to reflect one or more blocks
 allocated to files; and

20 establishing a third inode within said block map to said second block allocation
 bitmap subsequent to said modifying said second block allocation bitmap.

35. The method as recited in claim 28 wherein said establishing said second inode is
performed in response to a commit command.

25